

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problems Mailbox.**

⑩ 日本国特許庁(JP)

⑪ 特許出願公開

⑫ 公開特許公報(A) 平2-158871

⑬ Int.Cl.⁵

G 06 F 15/40

識別記号

5 0 0 T

庁内整理番号

7313-5B

⑭ 公開 平成2年(1990)6月19日

審査請求 未請求 請求項の数 2 (全4頁)

⑮ 発明の名称 文書分類装置

⑯ 特 願 昭63-312107

⑰ 出 願 昭63(1988)12月12日

⑱ 発 明 者 森 田 哲 也 東京都大田区中馬込1丁目3番6号 株式会社リコー内

⑲ 出 願 人 株 式 会 社 リ コ ー 東京都大田区中馬込1丁目3番6号

⑳ 代 理 人 弁 理 士 香 取 孝 雄 外1名

明 細 書

1. 発明の名称

文書分類装置

2. 特許請求の範囲

1. 文書データベースにおけるキーワードの出現頻度値を用いて計算される各キーワードの自己情報量を保持するキーワード情報量記憶手段と、

前記キーワードの自己情報量を用いて各文書ごとの概念特徴量を求める概念特徴抽出手段と、

文書間の該概念特徴量の差に応じて文書間の距離を求める文書間距離計算手段とを有し、

該文書間距離計算手段は、前記文書間の距離によって文書の分類を行なうことを特徴とする文書分類装置。

2. 文書データベースにおいて使用されるソーラスのキーワード分類項目ごとにキーワードの出現頻度値を用いて計算されるキーワードの自己情報量を保持するキーワード情報量記憶手段と、

各キーワード分類項目ごとの該キーワード情報量の総和をベクトル化したものを概念特徴量とし

て求める概念特徴抽出手段と、

文書間の該概念特徴量の差に応じて文書間の距離を求める文書間距離計算手段とを有し、

該文書間距離計算は、前記文書間の距離によって文書の分類を行なうことを特徴とする文書分類装置。

3. 発明の詳細な説明

〔産業上の利用分野〕

本発明は文書分類装置、とくに、文書に含まれるキーワードに基づき文書の概念特徴量を求め、概念特徴量により文書を分類する文書分類装置に関する。

〔従来の技術〕

文書をあらかじめ設定した分野へ自動的に分類するためカイ自乗値を用いてキーワードの偏りを調べ、文書を分類する方式が知られている。このような分類方式を記載したものとして、田村他「統計的手法による文書自動分類」(情報処理18回全国大会論文集、1987年)、および林知己夫「数量化の方法」(東洋経済新聞社、1974年)が

ある。

カイ自乗検定はキーワードの出現頻度の分野による偏りを示す指標としてカイ自乗値を求め文書を分類するものである。カイ自乗値は、各キーワードの出現頻度値と各分野ごとの総キーワード数が独立事象であると仮定した場合のキーワードの出現頻度値を理論度数とし、実測値との差を求め正規化したものである。

上記の文獻①はカイ自乗検定を用いて文書をあらかじめ設定した分野へ自動的に分類する方式について述べたものである。この方式は、キーワードの出現頻度の偏りを用いるために、あらかじめ大量の標本データを分野別に分類してカイ自乗値を計算し、分類用データを用意しておく必要がある。

文獻②もやはりカイ自乗値を用いる統計的手法の一つであり、複数の分野間の相関を見るための方式である。

〔発明が解決しようとする課題〕

上記の文獻①②に記載された方式は、標本デー

〔作 用〕

本発明によれば、キーワード情報量記憶手段が文書データベース等のキーワード出現頻度により、所定の計算を行って各キーワードの自己情報量を求め、概念特徴抽出手段が自己情報量より所定の計算により各文書の概念特徴量を求め、文書間距離計算手段が概念特徴量の差に応じて文書の分類を行なう。以上のようにキーワードの頻度より各手段の計算処理を通して、自動的に文書が分類されるので、従来の人手作業が不要となり、ばらつきのない、概念量による文書分類が構築できる。

〔実施例〕

本発明の実施例を図面を用いて具体的に説明する。

本発明による文書分類装置の一定実施例が図に示されている。

キーワード情報量記憶部1は入力される未登録文書Qよりキーワードを抽出し、後述のようにその出現頻度よりキーワードの出現確率を求め、そ

の分類にはやはり人手による作業が必要となる。したがって、人手による分類のばらつきや不適切さが介入するという問題がある。

また、後者は分類用の軸を決定するのが難しいという問題がある。

本発明は上記の問題点を解決するために、文書に含まれるキーワードの頻度値から各文書の概念特徴量を求め、これに応じて文書を分類する文書分類装置を提供することを目的とする。

〔課題を解決するための手段〕

上記目的を達成するために、本発明によれば、文書データベースにおけるキーワードの出現頻度値を用いて計算される各キーワードの自己情報量を保持するキーワード情報量記憶手段と、キーワードの自己情報量を用いて各文書ごとの概念特徴量を求める概念特徴抽出手段と、文書間の概念特徴量の差に応じて文書間の距離を求める文書間距離計算手段とを有する。

文書間距離計算手段は、文書間の距離によって文書の分類を行う。

の対数値をキーワード情報量Iとして記憶する。概念特徴抽出部2はキーワード情報量記憶部1よりキーワード情報量Iを入力し、その総和を文書Qの概念特徴量C(q)として出力する。文書間距離計算部3は概念特徴抽出部2より各文書の概念特徴量C(q)を入力して記憶し、2つの文書間の概念距離を求めて、概念距離の近い文書をクラスタ(分類)して、各種の分類を文書データベース4に格納する。各機能部は、各部の生成したデータを転送するデータバスa~cによって接続されている。

一般にシソーラス等のキーワード集に登録されているキーワードは、それらが現われる文書数や全文書における延べ出現回数等によって各キーワードの出現頻度を定義できる。いまキーワードKEYiの出現頻度を全キーワード数で正規化した P_i をキーワードKEYiの出現確率とすると、キーワードに出現確率 P_i を対応させるシステムは完全事象系となり以下のように表せる。

$$\begin{bmatrix} \text{KEY}_1 & \text{KEY}_2 & \dots & \text{KEY}_i & \dots & \text{KEY}_j & \dots & \text{KEY}_n \\ P_1 & P_2 & & P_i & & P_j & & P_n \end{bmatrix}$$

ただし、 $\sum_{i=1}^n P_i = 1$ である。

ここで、KEY_iの自己情報量 $I(\text{KEY}_i)$ は次式で表せる。

$$I(\text{KEY}_i) = -\log P_i \quad \dots (1)$$

また自己情報量は加法性を保つため、KEY_iとKEY_jの持つ合成情報量は、次式で表わされる。

$$\begin{aligned} I(\text{KEY}_i, \text{KEY}_j) &= I(\text{KEY}_i) + I(\text{KEY}_j) \\ &= -\log P_i - \log P_j \quad \dots (2) \end{aligned}$$

キーワード情報量記憶部1は、文書データベース4への未登録文書Qを概念特徴抽出部2を介してデータベースaより入力し、文書Qの各キーワードを抽出し、その出現確率KEY_iを求め、(1)式によりキーワードの自己情報量 $I(\text{KEY}_i)$ を計算して保持する。シソーラスが用意されているときは、シソーラスのキーワード分類項目ごとにキーワードの出現確率を求め、(1)式により自己情報量を計

2は、(3)式または(4)式を用いて、概念特徴量 $C(q)$ または $CV_r(q)$ を計算し、データベースbより文書間距離計算部3に出力する。

(3)式によって求められた概念情報量はある文書のもつキーワード情報量の和であり、その文書に付加された自己情報量の大きさを示しているだけである。この場合の概念情報量は、文書データベースの検索時における当該文書の分離度の高さ(測定しやすさ)を表す。このような分離度の高さによって文書を分類することも可能である。

しかし、通常は文書の内容によって既存の分類項目等に分類する用途が考えられる。そのような場合、(4)式の概念特徴量ベクトルを用いる。一般にM個の分類項目によってデータベースはM次元の概念空間を構成すると考えられる。従ってこのようなデータベース中の文書の持つ概念は、M個の特徴パラメータからなるM次元ベクトルとして表現できる。また任意の2つの概念特徴量ベクトルの距離が計算できるため、ある文書のある分類への帰属度や2つの文書間の概念的距離等が求

算できる。

ある文書Qのキーワード集合をqとしその概念特徴量を $C(q)$ と表すと、

$$C(q) = - \sum_{\text{KEY}_i \in q} \log P_i \quad \dots (3)$$

で与えられる。

また既存の分類項目を持つシソーラスにおいては概念特徴量をベクトルとして扱うことができる。最も単純な例として、M個の分類項目を持つシソーラスではM次元のベクトルCVを考える。今、R番目の分類項目に属するキーワードの集合をrとすると、文書Qの概念特徴量ベクトル $CV(q)$ のR要素 $CV_r(q)$ は、

$$CV_r(q) = - \sum_{i \in q \cap i \in r} \log P_i \quad \dots (4)$$

ただし、 $i \in q \cap i \in r$ はキーワードiが文書Q中に含まれ、かつR番目の分類項目中に含まれている場合の P_i の総和を計算することを意味する。

キーワード情報量記憶部1から文書Qの各キーワードの自己情報量Iを入力し、概念特徴抽出部

められる。

例えば、 $CV(q)$ という概念特徴量ベクトルを持つ文書が、キーワード集合kをもつ分類Kに帰属する度合を $INC(k, q)$ とすると、

$$INC(k, q) = CV_k(q) / \sum_{r=1}^M CV_r(q) \quad \dots (5)$$

で与えられる。

また、 $CV(s)$ 、 $CV(t)$ という概念特徴量ベクトルを持つ2つの文書間の概念距離を $D(s, t)$ とし例えば市街地距離で計算すると、

$$D(s, t) = \sum_{r=1}^M |CV_r(s) - CV_r(t)| \quad \dots (6)$$

で与えられる。

文書間距離計算部3は概念特徴量 $C(q)$ または $CV(q)$ を入力し、(5)式で示した計算を行なうことにより、未分類の文書の属すべき分類を決定でき、また(6)式を用いると、概念距離の近い文書群によっていくつかの分類を構成できる。文書間距離計算部3は文書Qの分類を文書データベース4に入力する。このとき生成される分類は、既存

のいくつかの分類項目の概念を結合した合成概念になるため、既存の分類項目に補われない文書概念自体に指向した新しい分類体系を自然に構築していく。

(5) 式を用いた同類文書の分類方法について具体的に説明する。

前述のように既存の分類項目に対して文書分類を行なう場合には、(5) 式を用いて各分類 K に属する度合い $INC(k, q)$ を求めればよい。さらに概念特徴量ベクトルを用いると、既存の分類項目を用いて新しい分類体系を構築することが可能となる。

まず、分類しようとする全ての文書について各文書間の概念距離 D を求める。次に全ての文書の中から任意に 1 文書 (文書 S とする) を選択し、その文書との概念距離が所定のしきい値より小さい、すなわちその文書と概念的に近い文書を抽出する。抽出された文書 T の集合を式で表現すると、文書 S 、 T に含まれるキーワード集合をそれぞれ、 s とすれば、

$$(T | D(s, t) < \theta)$$

(ただし、 $D(s, s) = 0$ は (5) より明らかであり、文書 S は必ず集合 T に含まれる。)

この作業を全ての文書に対して行なうと文書数に等しい同類文書の集合が出来上がる。これら同類文書集合をその集合の要素数 (文書数) に従って降順に並べ、文書数の多い順に必要な分類数だけの同類文書集合を選択する。この選択は分類数で制限しても良いし、文書数で制限しても良い。分類可能な数の最大値は文書数である。この場合各分類に含まれる文書数は 1 であるが、このような分類が最適となる場合もあってしかるべきである。

【発明の効果】

本発明によれば、キーワード抽出、または既存のキーワード集の分類を用いて概念特徴量を計算できるため、未登録文書の分類の前に評価用データを作成する必要がない。

概念距離の近い文書群によって分類を構成する

ため、既存の分類項目に補われない文書概念自体に指向した新しい分類体系を自然に構築していくという優れた効果がある。

4. 図面の簡単な説明

図は本発明の文書分類装置の一実施例を示す機能ブロック図である。

主要部分の符号の説明

- 1…キーワード情報量記憶部、
- 2…概念特徴抽出部、
- 3…文書間距離計算部、
- 4…文書データベース。

特許出願人 株式会社リコー

代理人 香取 孝雄
丸山 隆夫

